

1. [Introduction](#)
2. [Project Overview](#)
3. [Background on Independent Component Analysis](#)
4. [Blind Source Separation](#)
5. [Background on Artificial Neural Network](#)
6. [Implementation and Experimental Results](#)
7. [Conclusion and Next Steps](#)
8. [Team Members and Acknowledgements](#)
9. [References](#)
10. [Poster](#)

Introduction

Description of the motivation for undertaking the task of creating selective transparent headphones.

The advancement of signal processing technologies has consistently improved the enjoyment of listening to music. For example, noise cancelling techniques have been widely applied to headphones that, in addition to physical noise isolation, provide further reduction in noise.

We believe, however, that it is sometimes important for music listeners with a pair of good noise isolating headphone to be aware of certain surrounding signals. For example, when people are crossing the street while listening to music, they really want to hear car horns to be warned of possible danger. We also think that human speech is important, and that people may want to hear others talking to them while they are listening to music.

Based on the assumption above, we propose to build a selective transparent headphone that propagates certain surrounding sound signals and suppresses all other signals. In this project, we explore possibility to select only speech signal to propagate and silence other surrounding sounds.

The following sections are organized as follows. The second section provides an overview of the project, defining the problem and explaining high level system implementation. The third to the fifth sections explain in detail the blind source separation algorithm and binary artificial neural network, respectively. The following sections discuss experiment results, and suggest next steps.

Project Overview

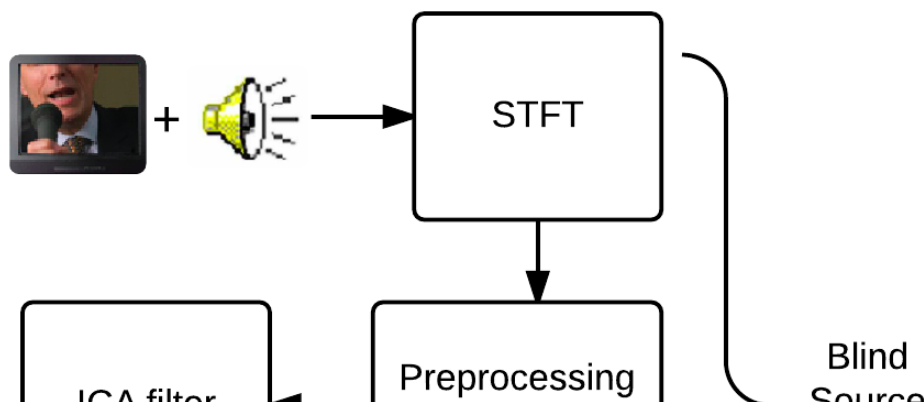
Overview of the steps followed to create the system.

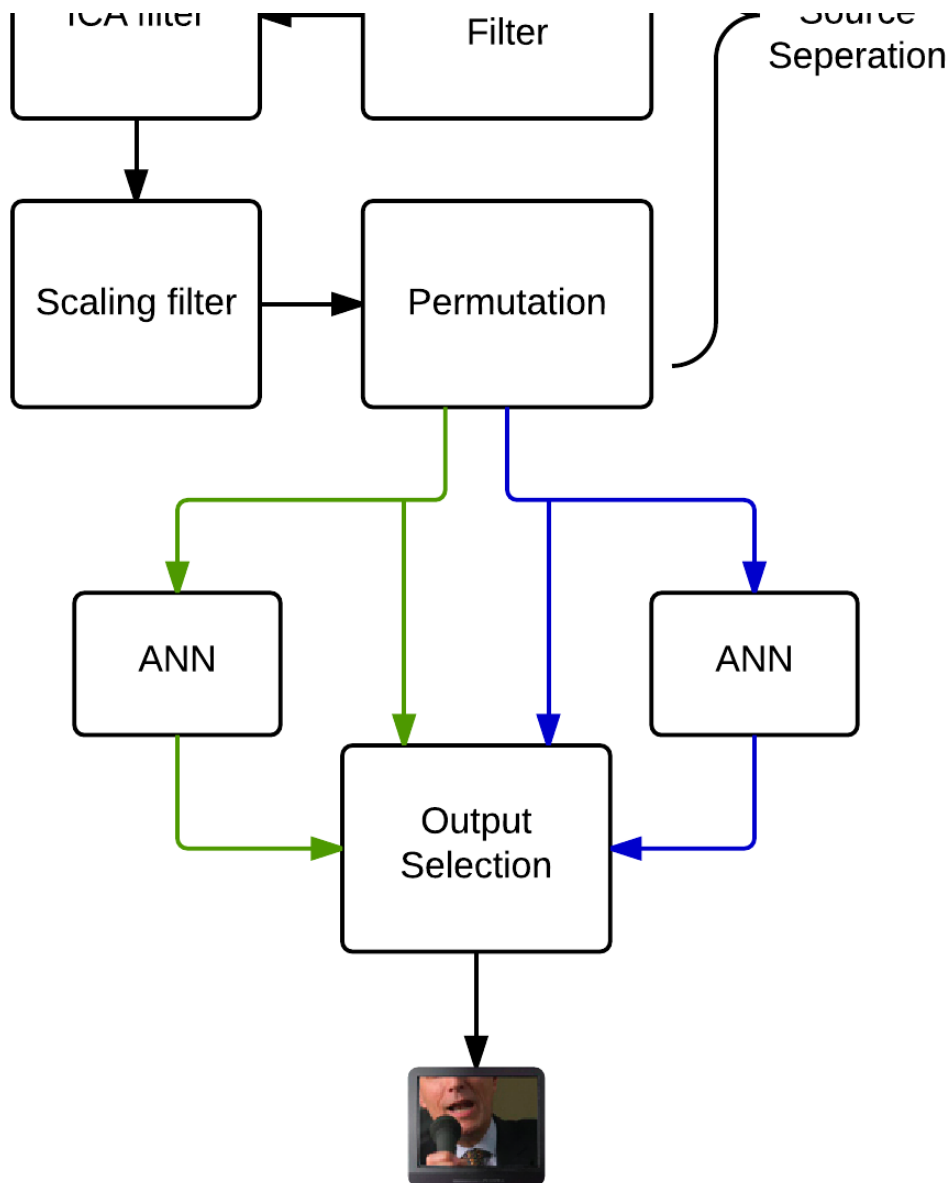
Problem Definition

There are a number of sound signals that are of importance in the environment. In this project, we identify human speech as an important signal, that we would like to focus on. We would like to separate the environment sound signals into a signal containing only human speech and a signal containing all other speech. By making the above simplifications, the problem is reduced to finding the speech content in the surrounding environment, and forwarding the speech content to the listener while suppressing other signals. If there are no speech signals, or speech signals are weak compared to other signals, all sounds from the environment will be attenuated.

System Implementation

We approach the problem stated above by first separating the source signals into human speech content and non-human speech with a blind source separation algorithm. After that, a classification algorithm determines which signal contains human speech, if any, and outputs that signal. A detailed overview is shown below with the system block diagram in Figure 1.





The audio input, a mixture of human speech and some other noise such as instrumental music, is passed to the blind source separation block. Inside this block, short time Fourier transform is first performed to change the time domain input signal into frequency domain, since all the other operations on the signals reside in frequency domain. After that a preprocessing filter is applied, which cleans the input signal by removing reflection and ambient noises. Then for each frequency, the independent component analysis algorithm separates the signal into two parts such that independence between the two signals is maximized. Scaling and

permutation filters minimizes distortions caused by using the independent component analysis method. The blind source separation process outputs two signals, where at most, one signal contains speech.

We then implement a binary artificial neural network (ANN) for classification. The two identical ANN take the two signals separated from the source as inputs respectively and outputs a weight for each signal. We train the neural network in the way that a signal containing more human speech will output a higher weight. Based on the weights of the separated signals, the output selection multiplexer chooses the proper signal to output. That is, if one signal, referred to as A, has a higher weight than the other signal, referred to as B, and its weight is above a certain threshold (0.5), the multiplexer outputs signal A. If both the weights of signal A and signal B are below that threshold, then the multiplexer outputs nothing. We set a threshold to deal with situation where both signal contain no human speech.

Background on Independent Component Analysis

Describes the motivation and background behind using independent component analysis as the basis of the blind source separation algorithm used in this endeavor.

The blind source separation algorithm, we employed, is based around independent component analysis, which is explained in more detail below. The goal is to recover independent sources given only sensor observations that are linear mixtures of independent source signals. We assume that the source signals are statistically independent and non-Gaussian.

ICA works by solving the following model.

$$\mathbf{x} = \mathbf{A}\mathbf{s}$$

Where vector \mathbf{x} contains n observed signals, vector \mathbf{s} contains independent sources that comprise \mathbf{x} , and matrix \mathbf{A} denotes the mixing matrix. It is assumed that the mixing matrix \mathbf{A} and the source vector \mathbf{x} are both unknown. ICA algorithm thus produces the best estimation for both the mixing matrix and source vector.

The first step in ICA is usually whitening (sphering) the data, which removes any correlations in the data, i.e. the signals are forced to be uncorrelated. Putting the words in mathematical terms, we seek a linear transformation \mathbf{V} such that when $\mathbf{y} = \mathbf{V}\mathbf{x}$ we now have $\mathbf{E}[\mathbf{y}\mathbf{y}'] = \mathbf{I}$.

After sphering, the separated signals can be found by an orthogonal transformation of the whitened signals \mathbf{y} (this is simply a rotation of the joint density). The appropriate rotation is sought by maximizing the non-normality of the marginal densities. This is because of the fact that a linear mixture of independent random variables is necessarily more Gaussian than the original variables.

Although in theory ICA algorithm can perfectly separate source signals from observed signals, provided that sources are statistically independent and non-Gaussian, most of the algorithms do not work quite well in real world scenarios. For example, ICA does not typically deal with reflection; that is, if mixed signal is observed in a small room with reverberation, ICA algorithms such as FastICA hardly identifies individual components from

the observed signal. Therefore, we need a more robust blind source separation algorithm, which is discussed in the following section.

Blind Source Separation

Describes the method to separate an arbitrary signal into its independent components for the application of creating a selective transparent headphone.

Background

When separating mixtures of instantaneously mixed signals, independent component analysis works very well, but this ideal situation is rarely found in the real-world. In practical applications, the environment distorts audio signals by adding echoes, reflections, and ambient noise. Additionally independent component analysis, in its purest form, assumes that source signals do not have any propagation delay, which is an assumption that cannot be applied in this case. Recording sources from two microphones placed in different locations will inevitably introduce propagation delays, so the blind source separation method used must also consider this issue.

To solve problems detailed above, the blind source separation problem will be redefined in the time-frequency domain. By taking the short-time Fourier transform (STFT) of the audio inputs, we can represent the inputs as the following.

$$x(\omega, t) = [X_1(\omega, t), X_2(\omega, t), \dots, X_M(\omega, t)]^T$$

$\mathbf{X}(\omega, t)$ refers to the input observed at the i th microphone observed at time t . The T symbol refers to the transpose, so in this case, the input sources are represented along the rows. For our application, we will assume that the number of observed inputs and the number of separated sources are both equal to the constant M . We can, now, formulate our blind source separation problem as the following.

$$x(\omega, t) = A(\omega)s(\omega, t) + n(\omega, t)$$

In this case, $\mathbf{s}(\omega, \mathbf{t})$, refers to the vector of source signals observed at frequency ω and at time t . The second term, $\mathbf{n}(\omega, \mathbf{t})$, represents any type of noise or distortions that may be present in the observed signals (noise, reflections, etc.). The blind source separation method attempts to solve for the mixing matrix, $\mathbf{A}(\omega)$, which can be represented as the following.

$$A_{i,j}(\omega) = H_{i,j}(\omega)e^{-j\omega\tau_{i,j}}$$

Here, $H_{i,j}(\omega)$, refers to the transfer function of the j th source to the i th microphone. Additionally, $\tau_{i,j}$, refers to the propagation delay from the j th source to the i th microphone.

The system, for which our problem is defined on, is now redefined from an instantaneous mixture of signals to a convolutional mixture in the time-domain of the following form.

$$x(t) = \hat{a}(t) * s(t)$$

This formulation allows for a solution which considers the distortion, $\mathbf{n}(\omega, \mathbf{t})$, added by the environment and considers the inherent propagation delay that is inherently present in our application.

Since the new formulation of the blind source separation problem is a convolutional system in the time-domain, we will solve the separation problem in the time-frequency domain formulation shown in Figure 2. Given our approach outlined above, we can solve for the complete frequency-domain separation filter shown below.

$$G(\omega) = P(\omega)\tilde{B}_m^+(\omega)B(\omega)$$

$\mathbf{B}(\omega)$ is the separation filter of the system. Applying this filter will separate the spectrum of the observed signal into the spectra of the source signals

observed at frequency, ω . As stated above, by reformulating the problem into the frequency domain, the convolutional mixture becomes transformed into an instantaneous mixture. As such, independent component analysis can be applied to separate the source signals at each frequency. However, in doing so, we introduce other problems, inherent to the independent component analysis method, that the other filters attempt to resolve.

The following sections will explain the three filters in detail.

- **B(ω)**, the separation filter.
- **Bm+**, the normalization filter used to solve the scaling problem that arises after applying independent component analysis.
- **P(ω)**, the permutation filter used to solve the frequency distortion problem that arises after applying independent component analysis.

Separation Filter

The separation filter, shown in the previous section can be broken down into the following components.

$$B(\omega) = U(\omega)W(\omega)$$

W(ω) is the preprocessing filter used to reduce reflections and ambient noise by using the subspace method. **U(ω)** is the filter obtained by applying independent component analysis after preprocessing.

I. Subspace Method

The spatial correlation, or autocorrelation, matrix at frequency ω is defined below.

$$R(\omega) = E[x(\omega, t)x^H(\omega, t)]$$

Given that there are M inputs and M sources, the resulting matrix will be $M \times M$. This matrix can be rewritten in the following way.

$$R(\omega) = AQA^H + K$$

The spatial
correlation matrix.

$$K = E[n(t)n^H(t)]$$

The noise
correlation
matrix.

$$Q = E[s(t)s^H(t)]$$

The sources'
cross-spectrum
matrix.

In other words, \mathbf{K} is the correlation matrix of $\mathbf{n}(\mathbf{t})$ and \mathbf{Q} is the cross-spectrum matrix of the sources. However, since the source signals, $\mathbf{s}(\mathbf{t})$, are unknown at this point, this equation cannot be directly solved. Instead, we can represent the generalized eigenvalue decomposition of the spatial correlation matrix in the following manner.

$$\mathbf{R} = \mathbf{K} \mathbf{E} \mathbf{\Lambda} \mathbf{E}^{-1}$$

\mathbf{E} refers is the eigenvector matrix, $\mathbf{E} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_M]$, and $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_M)$ refers to the eigenvalues of \mathbf{R} . Since \mathbf{K} , the noise correlation matrix, cannot be directly observed apart from the source signals, we will assume that $\mathbf{K} = \mathbf{I}$, which will evenly distribute the reflections and ambient noise induced by the environment among the estimated sources. As such, the eigenvalue decomposition of \mathbf{R} can be rewritten as the following.

$$\mathbf{R} = \mathbf{E} \mathbf{\Lambda} \mathbf{E}^{-1}$$

The final preprocessing filter uses the eigenvector and eigenvalue matrices resulting from the standard eigenvalue decomposition of the spatial correlation matrix and is shown below.

$$\mathbf{W} = \mathbf{\Lambda}^{-1/2} \mathbf{E}^H$$

II. Independent Component Analysis

By applying the preprocessing filter obtained from the subspace method, our observed signals are now of the following form.

$$y(\omega, t) = \mathbf{W}(\omega) x(\omega, t)$$

To estimate the source signals from the preprocessed input signals, we can use independent component analysis to solve the linear system displayed below, for the filter $\mathbf{U}(\omega)$.

$$y(\omega, t) = U^{-1}(\omega)\hat{s}(\omega, t)$$

Now that the separation filter at each frequency has been found, the problems that arise from the separation process must be addressed.

Scaling

The scaling problem that results from applying independent component analysis can be resolved using the filter shown below.

$$\tilde{B}_m^+(\omega) = \text{diag}[B_{m,1}^+, B_{m,2}^+, \dots, B_{m,M}^+]$$

B_{m,n}⁺(ω) refers to the (m,n) th entry in **B⁺(ω)**, which is the pseudo-inverse of **B(ω)** (regular inverse also works here since we assume that the number of separated sources and the number of inputs are both M , so that the resulting separation matrix is square). Additionally, m refers to an arbitrary row or microphone in **B⁺(ω)**. By applying this matrix, each signal, i , is amplified by the component of signal i observed at microphone or input m . Essentially, this filter amplifies the frequency components of the separated source signals so that the waveforms of the resulting sources will be distinguishable and audible.

Permutation

A second critical problem that arises after applying independent component analysis is that the order in which the source signals are returned is unknown. Since independent component analysis is applied at each frequency, we must find the permutation of components that has the highest chance of being the correct permutation to reconstruct the source signals and minimize the amount of frequency distortion caused by independent component analysis.

We define the permutation matrix, **P**, as the following.

$$Z(\omega) = B^+(\omega)$$

$$\bar{Z}^T(\omega) = PZ^T(\omega)$$

$$\bar{Z}(\omega) = [\bar{z}_1(\omega), \bar{z}_2(\omega), \dots, \bar{z}_M(\omega)]$$

The matrix, \mathbf{P} , exchanges the column vectors of $\mathbf{Z}(\omega)$ to get different permutations. We define the cosine of θ_n between the two vectors $\mathbf{z}_n(\omega)$ and $\mathbf{z}_n(\omega_0)$, where ω_0 is the reference frequency is as follows.

$$\cos(\theta_n) = \frac{\bar{z}_n^H(\omega) \bar{z}_n(\omega_0)}{\|\bar{z}_n(\omega)\| \cdot \|\bar{z}_n^H(\omega_0)\|}$$

The permutation is, then, determined by the following.

$$P = \arg \max_P F(P)$$

The cost function, $\mathbf{F}(\mathbf{P})$, used above, is defined as follows.

$$F(P) = \frac{1}{D} \sum_{n=1}^M \cos(\theta_n)$$

This compares the inverse vector filter at each frequency with a filter at a reference frequency, ω_0 . However, if we use the same reference frequency to find the permutation at every other frequency, then the problem arises if the filter at the reference frequency has the incorrect permutation. In order to minimize this error, a new reference frequency is chosen after the permutations of filters at \mathbf{K} frequencies are found. As such, the frequency range of a reference frequency, ω_0 , is as follows.

$$\omega_0 = \omega - k \cdot \Delta\omega, \forall k = 1, \dots, K$$

Background on Artificial Neural Network

This module talks about the backgrounds on artificial neural network

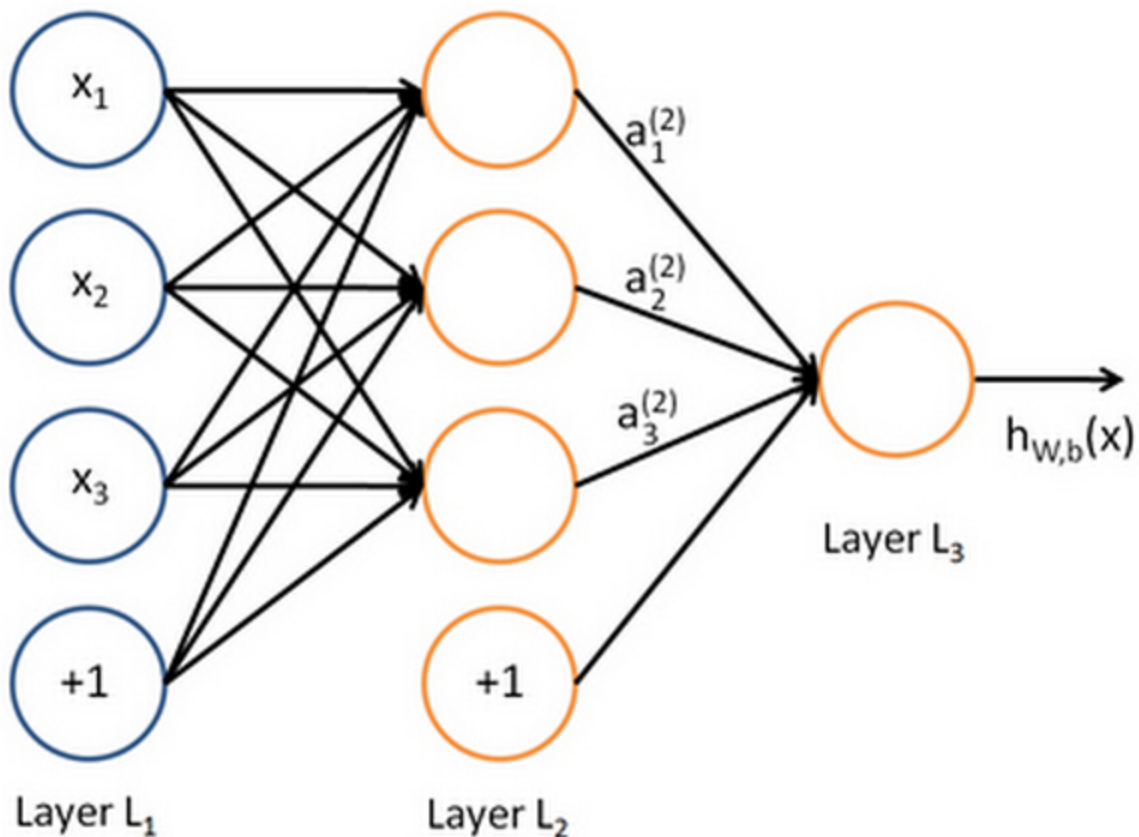
Motivation:

As the Blind Source Separation algorithm randomly assigns the two independent components into two output vectors, we need the artificial neural network to determine which of the two signals, if any, contains human speech.

Background on Our Chosen Neural Network Model:

Our Neural Networks Model:

A simplified schematics of our neural network looks like the following:



Neural Network Model

Our artificial neural network consists of one input layer (Layer L1), one hidden layer (Layer L2), and one output layer (Layer L3). The circles labeled “+1” are called the bias units, and corresponds to the intercept term.

Our neural network has parameters $(W, b) = (W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)})$, where we write $W^{(l)}_{ij}$ to denote the weight associated with the connection between unit j in layer l , and unit i in layer $l+1$. Also, $b^{(l)}_i$ is the bias associated with unit i in layer $l+1$. We also use s_l to denote the number of nodes in layer l (not counting the bias unit).

We will write $a^{(l)}_i$ to denote the activation of unit i in layer l . For $l=1$, we also use $a^{(1)}_i = x_i$ to denote the i -th input. Given a fixed setting of the parameters W, b , our neural network defines a hypothesis $h_{W,b}(x)$ that outputs a real number. Specifically, the computation that this neural network represents is given by:

$$\begin{aligned} a^{(2)}_1 &= f(W^{(1)}_{11} x_1 + W^{(1)}_{12} x_2 + W^{(1)}_{13} x_3 + b^{(1)}_1) \\ a^{(2)}_2 &= f(W^{(1)}_{21} x_1 + W^{(1)}_{22} x_2 + W^{(1)}_{23} x_3 + b^{(1)}_2) \\ a^{(2)}_3 &= f(W^{(1)}_{31} x_1 + W^{(1)}_{32} x_2 + W^{(1)}_{33} x_3 + b^{(1)}_3) \\ h_{W,b}(x) &= a^{(3)}_1 = f(W^{(2)}_{11} a^{(2)}_1 + W^{(2)}_{12} a^{(2)}_2 + W^{(2)}_{13} a^{(2)}_3 + b^{(2)}_1) \end{aligned}$$

Let $z^{(l)}_i$ denote the total weighted sum of inputs to unit i in layer l , including the bias unit. For instance,

$$z^{(2)}_i = \sum_{j=1}^n W^{(1)}_{ij} x_j + b^{(1)}_i$$

, so that

$$a^{(l)}_i = f(z^{(l)}_i).$$

.

If we extend the activation function $f(\cdot)$ to apply to vectors in an element-wise fashion, we can write the above equations more compactly as:

$$\begin{aligned}
z^{(2)} &= W^{(1)}x + b^{(1)} \\
a^{(2)} &= f(z^{(2)}) \\
z^{(3)} &= W^{(2)}a^{(2)} + b^{(2)} \\
h_{W,b}(x) &= a^{(3)} = f(z^{(3)})
\end{aligned}$$

The above equations represents the forward propagation step. More generally, as we use $a^{(1)} = x$ to denote the values from the input layer, then given layer l 's activations $a^{(l)}$, we can compute layer $l+1$'s activations $a^{(l+1)}$ as:

$$\begin{aligned}
z^{(l+1)} &= W^{(l)}a^{(l)} + b^{(l)} \\
a^{(l+1)} &= f(z^{(l+1)})
\end{aligned}$$

Backpropagation Algorithm:

Suppose we have a fixed training set

$$\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$$

of m training examples, the overall cost function is:

$$\begin{aligned}
J(W, b) &= \left[\frac{1}{m} \sum_{i=1}^m J(W, b; x^{(i)}, y^{(i)}) \right] + \frac{\lambda}{2} \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (W_{ji}^{(l)})^2 \\
&= \left[\frac{1}{m} \sum_{i=1}^m \left(\frac{1}{2} \|h_{W,b}(x^{(i)}) - y^{(i)}\|^2 \right) \right] + \frac{\lambda}{2} \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (W_{ji}^{(l)})^2
\end{aligned}$$

The first term in the cost function is an average sum-of-squares error term. The second term is the regularization term which is used to prevent

overfitting. The weight decay parameter λ controls the relative importance of the cost term and the regularization term.

Our goal is to minimize the cost function $J(W,b)$ as a function of W and b . To begin training our neural network, we will initialize each weight $W^{(l)}_{ij}$ and each $b^{(l)}_i$ to a small value near zero. It is important to initialize the parameters randomly for the purpose of symmetry breaking.

In order to perform gradient descent to minimize the cost function, we need to compute the derivatives of the cost function with respect to $W^{(l)}_{ij}$ and $b^{(l)}_i$ respectively:

$$\frac{\partial}{\partial W^{(l)}_{ij}} J(W, b) = \left[\frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial W^{(l)}_{ij}} J(W, b; x^{(i)}, y^{(i)}) \right] + \lambda W^{(l)}_{ij}$$

$$\frac{\partial}{\partial b^{(l)}_i} J(W, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial b^{(l)}_i} J(W, b; x^{(i)}, y^{(i)})$$

The intuition behind the backpropagation algorithm is as follows. Given a training example (x,y) , we will first run a “forward propagation” to compute all the activations, including the output value of the hypothesis $h_{W,b}(x)$. Then for each node i in layer l , we compute an error term $\delta^{(l)}_i$ that measures how much that node was responsible for any errors in the output. For an output node, we can directly measure the difference between the network’s hypothesis and the true value, and use that to define the error term for the output layer.

Here is the implementation of the backpropagation algorithm in MatLab:

1. Perform a feedforward propagation, computing the activations for hidden layer L2 and output layer L3.
2. For the output layer (layer n_l), set

$$\delta^{(n_l)} = -(y - a^{(n_l)}) \bullet f'(z^{(n_l)})$$

1. For the hidden layer, set

$$\delta^{(l)} = ((W^{(l)})^T \delta^{(l+1)}) \bullet f'(z^{(l)})$$

1. Compute the desired partial derivative of the cost function with respect to W and b

$$\begin{aligned}\nabla_{W^{(l)}} J(W, b; x, y) &= \delta^{(l+1)} (a^{(l)})^T, \\ \nabla_{b^{(l)}} J(W, b; x, y) &= \delta^{(l+1)}.\end{aligned}$$

Implementation note: in step 2 and 3 above, we need to compute $f'(z^{(l)}_i)$ for each value of i. as our $f(z)$ is the sigmoid function and we already have $a^{(l)}_i$ computed during the feedforward propagation process. Thus using the expression for $f'(z)$, we can compute this as

$$f'(z^{(l)}_i) = a^{(l)}_i(1 - a^{(l)}_i).$$

.

Below is the pseudo code of the gradient descent algorithm:

Note: $\Delta W^{(l)}$ is a matrix of the same dimension as $W^{(l)}$ and $\Delta b^{(l)}$ is a vector of the same dimension as $b^{(l)}$. one iteration of the gradient descent as follows:

1. Set $\Delta W^{(l)} := 0$, $\Delta b^{(l)} := 0$ for all l.
2. For i from 1 to m,
 - a. Use backpropagation to compute

$$\begin{aligned}\nabla_{W^{(l)}} J(W, b; x, y) \text{ and} \\ \nabla_{b^{(l)}} J(W, b; x, y).\end{aligned}$$

- a. Set

$$\Delta W^{(l)} := \Delta W^{(l)} + \nabla_{W^{(l)}} J(W, b; x, y).$$

- b. Set

$$\Delta b^{(l)} := \Delta b^{(l)} + \nabla_{b^{(l)}} J(W, b; x, y).$$

1. Update the parameters:

$$\begin{aligned} W^{(l)} &= W^{(l)} - \alpha \left[\left(\frac{1}{m} \Delta W^{(l)} \right) + \lambda W^{(l)} \right] \\ b^{(l)} &= b^{(l)} - \alpha \left[\frac{1}{m} \Delta b^{(l)} \right] \end{aligned}$$

We can now repeat the gradient descent steps to reduce our cost function $J(W, b)$.

Implementation and Experimental Results

This is our implementation of neural network and Blind Source Separation and their respective results

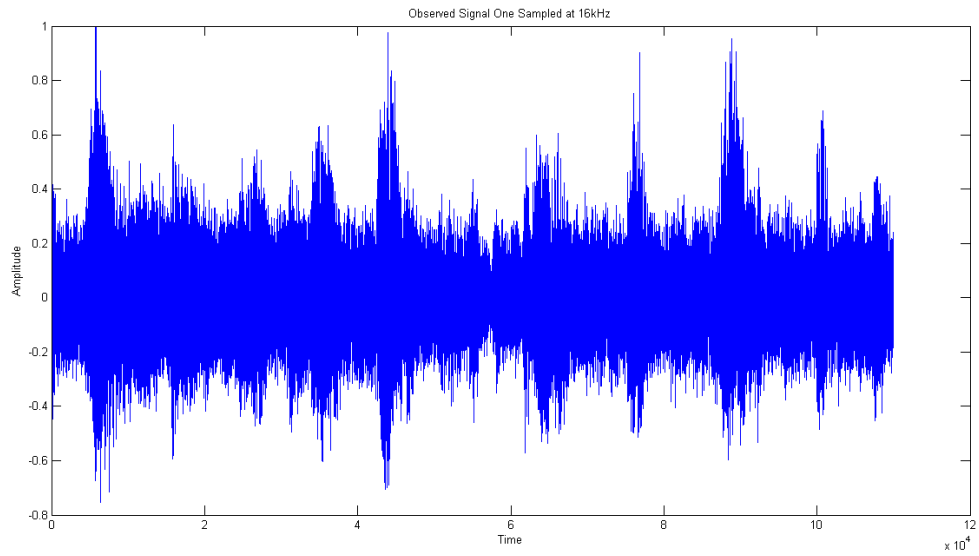
Implementation of Blind Source Separation:

When implementing the Blind Source Separation algorithm, we chose the following parameters,

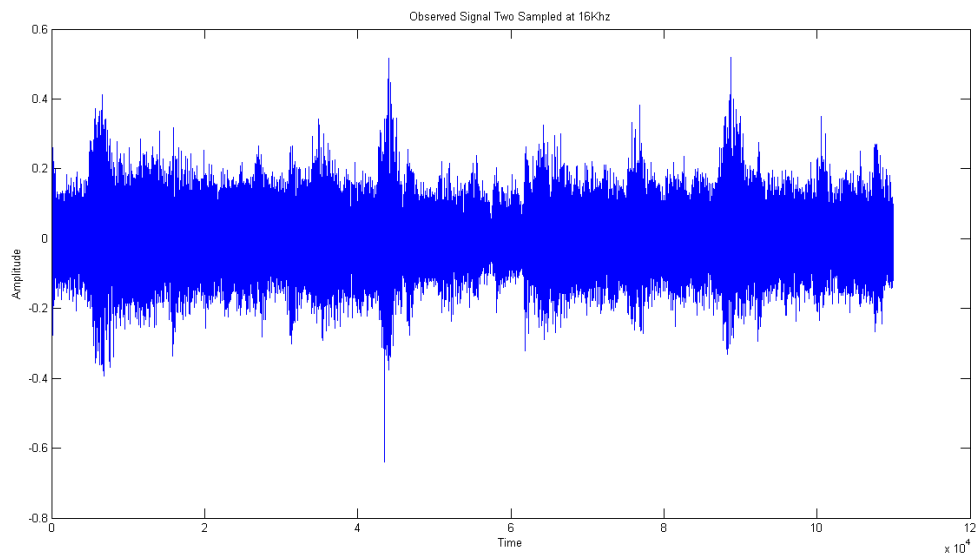
Parameters	
Sampling Rate	16 kHz
Length of STFT	512
Number of FFT Points	4096
Overlap	492
Number of inputs	2
Permutation	5
Reference Range	

Input Parameter Into The System

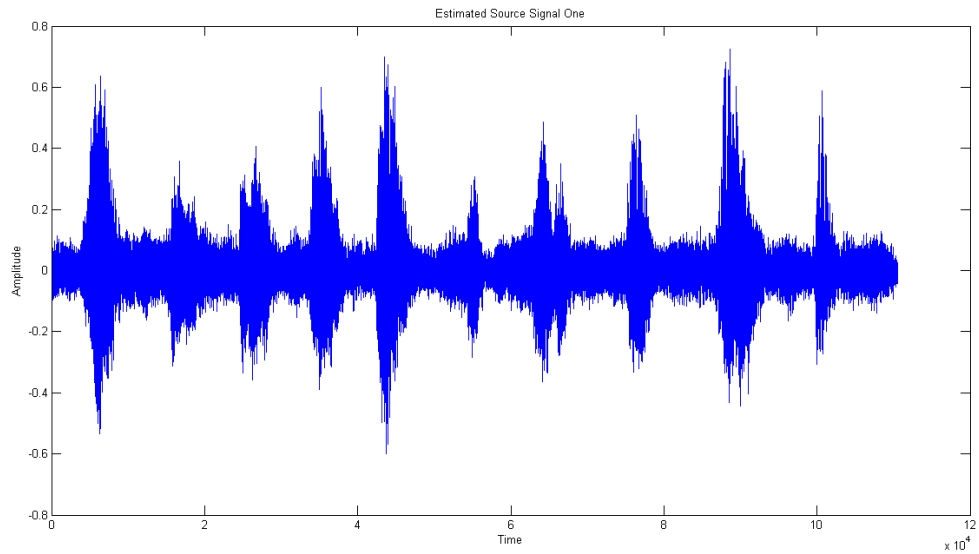
Result of our Blind Source Separation:



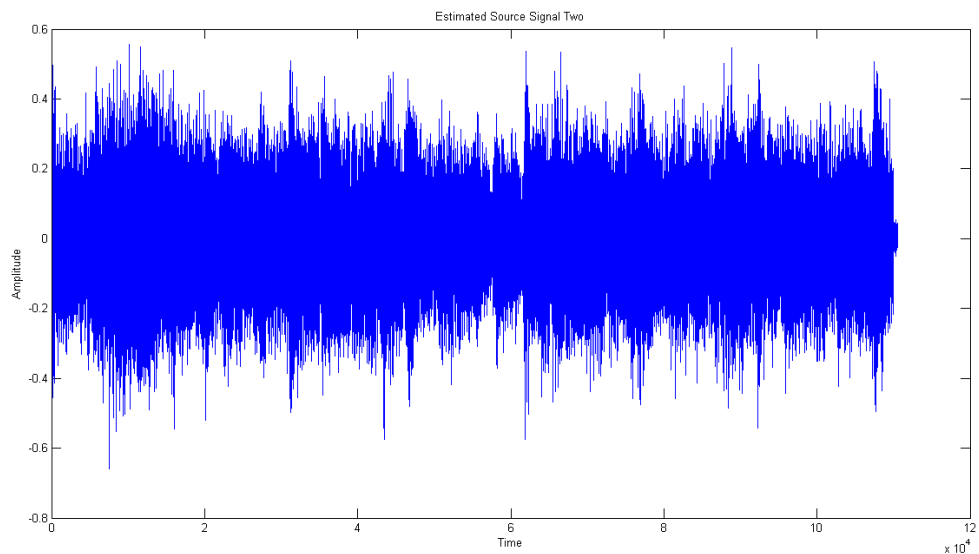
Mixed Signal Observed at 1st Microphone



Mixed Signal Observed at 2nd Microphone



Estimated Source Signal One-human speech



Estimated Source Signal Two-Instrumental Music

Figure 2 and Figure 3 show the mixed signals observed at the 1st microphone and the 2nd microphone respectively. The mixed signal is a

recording of a person counting from zero to ten with instrumental music playing in the background. The intermittent spikes in the mixed signal represents the person counting while the continuous, small-amplitude signal represents the instrumental music.

After applying the Blind Source Separation algorithm, we obtain two independent signals shown in Figure 4 and Figure 5 respectively. These two signals are very distinct from each other. The output signal 1 shown in Figure 4 consists mainly of the person counting whereas the output signal 2 shown in Figure 5 is primarily the instrumental music.

Implementation of Neural Network:

Since we need our neural network to perform binary classification, we used 512 nodes for the input layer, 25 nodes for the hidden layer, and 2 nodes for the output layer.

We used 10,000 positive training examples and 10,000 negative training examples to train our neural network. The positive training examples mainly consist of audio file containing examples of human speech such as news broadcasting. The negative training examples mainly consist of examples of non-human speech such as instrumental music.

After taking 512 point Short-Time Fourier Transform, we place our training sets into a 20,000-by-512 matrix where each row represents a single training example. Then we used the backpropagation and gradient descent algorithm described in the section above to train our neural network. We used 70% of the training examples for training, 15% of the training examples for validation, and another 15% for testing.

Result of Our Training of Neural Network:

After employing backpropagation algorithm and gradient descent algorithm to train our neural network, we get the following test results:

Training Confusion Matrix			
Output Class	1	2	
	11959 51.8%	190 0.8%	98.4% 1.6%
	2 0.0%	10949 47.4%	100.0% 0.0%
	1	2	
Target Class	100.0% 0.0%	98.3% 1.7%	99.2% 0.8%

Validation Confusion Matrix			
Output Class	1	2	
	2500 50.5%	108 2.2%	95.9% 4.1%
	28 0.6%	2314 46.7%	98.8% 1.2%
	1	2	
Target Class	98.9% 1.1%	95.5% 4.5%	97.3% 2.7%

Test Confusion Matrix			
Output Class	1	2	
	2476 50.0%	104 2.1%	96.0% 4.0%
	35 0.7%	2335 47.2%	98.5% 1.5%
	1	2	
Target Class	98.6% 1.4%	95.7% 4.3%	97.2% 2.8%

All Confusion Matrix			
Output Class	1	2	
	16935 51.3%	402 1.2%	97.7% 2.3%
	65 0.2%	15598 47.3%	99.6% 0.4%
	1	2	
Target Class	99.6% 0.4%	97.5% 2.5%	98.6% 1.4%

Accuracy of Neural Network

As our neural network is a binary classification model, our output is a vector containing two entries. The first entry represents the probability of the input signal being a human speech and the second entry represents the probability of the input signal being a non-human speech. The two entries sum to one. The overall performance of our neural network is shown in the

All Confusion Matrix. The red block showing 0.2% means 0.2% of the class 1 examples (human speech) has been incorrectly identified as class 1 (non-human speech). Similarly, the red block showing 1.2% mean 1.2% of the examples of class 2 (non-human speech) has been incorrectly identified as class 1 (human speech). The blue block shows that our neural network is able to distinguish between human speech and non-human speech with an accuracy of 98.6%.

Conclusion and Next Steps

This is our conclusion

Conclusion

From our experimental results, we can see that our Blind Source Separation algorithm is able to separate the two independent components of the mixed signal very well. After the separation, the neural network can distinguish between human speech and non-human sound with an accuracy of 98.6%. Consequently, our entire system can successfully output the human speech from a linear mixture of sound signals.

Next Steps

Implement the system in real-time. Our design currently works in the Matlab environment. It is implemented to deal with audio that is already recorded, but it does not deal with real-time audio streams. We eventually would like the system to work in real-time, separating the sources, forward speech, if any, and suppress other signals instantly so that people can benefit from this system. This implies improving the efficiency of the algorithms, and making use of Matlab SimuLink to explore the possibility of real-time implementation of our system.

Additionally, we would like to explore the possibility to separate several sources in addition to human speech. As mentioned in the introduction, other signals from the surrounding might also be important, such as a car's horn. We can improve the separation algorithm and train the neural network so that the system can recognize not only human speech, but also other important signals from the environment. With this implementation, users can have the option to forward in the signals they want. In order to achieve that, one needs to increase the number of microphones to increase the number of observed signals, and re-train the neural network so that it recognizes different types of sound signals.

Team Members and Acknowledgements

Team Members

Contact Information

Wang, Zichao: zichao.wang@rice.edu

Xia, Stephen: stephen.xia@rice.edu

Yao, Tianyi: tianyi.yao@rice.edu

Zhao, Shengjia: shengjia.zhao@rice.edu

Work Distribution

Blind Source Separation: Zichao and Stephen

Neural Network: Tianyi and Shengjia

Poster and Presentation: Zichao, Stephen, Tianyi, Shengjia

Documentation: Zichao, Stephen, Tianyi, Shengjia

Acknowledgements

We would like to thank Dr. Mohammad Golbabaee for mentoring and providing us with many insights throughout our project. Additionally, we would like to extend our thanks to Professor Richard Baraniuk and the Rice University Electrical and Computer Engineering Department for giving us the opportunity to partake in this endeavor.

References

Asano, Futoshi, et al. "Combined approach of array processing and independent component analysis for blind separation of acoustic signals." *Speech and Audio Processing, IEEE Transactions on* 11.3 (2003): 204-215.

"Backpropagation Algorithm." *Unsupervised Feature Learning and Deep Learning*. Stanford University, 26 Feb. 2011. Web. 16 Dec. 2014.

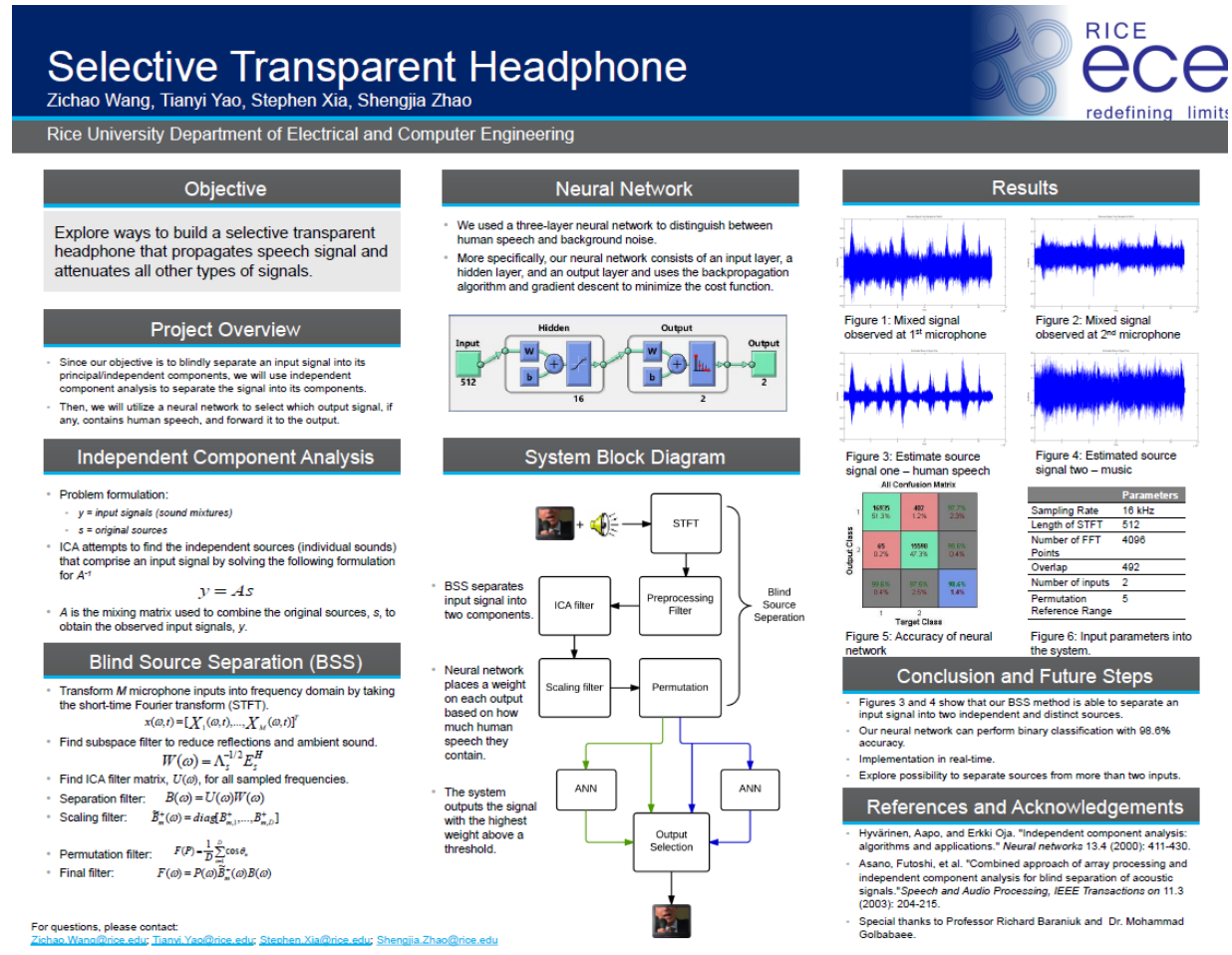
Hyvärinen, Aapo, and Erkki Oja. "Independent component analysis: algorithms and applications." *Neural networks* 13.4 (2000): 411-430.

"Neural Networks." *Unsupervised Feature Learning and Deep Learning*. Stanford University, 26 Feb. 2011. Web. 16 Dec. 2014.

Widrow, Bernard, et al. "Adaptive noise cancelling: Principles and applications." *Proceedings of the IEEE* 63.12 (1975): 1692-1716.

Poster

This is our poster



Our Project Poster